

INTRODUCTION TO MACHINE LEARNING COMPSCI 4ML3

LECTURE 11

HASSAN ASHTIANI

MAXIMUM LIKELIHOOD ESTIMATE

- MAXIMIZES THE PROBABILITY OF THE OBSERVATIONS GIVEN THE PARAMETERS

- $\alpha^{ML} = \underset{\alpha}{\operatorname{argmax}} P(X|\alpha)$

- $\alpha^{ML} = \underset{\alpha}{\operatorname{argmin}} -\left(\sum_i \operatorname{LOG} P(x^i|\alpha)\right)$

- FOR BIAS OF THE COIN

- $\alpha^{ML} = \frac{n_0}{n_0+n_1}$

MAXIMUM A POSTERIORI ESTIMATE

- MAXIMIZES THE PROBABILITY OF THE PARAMETERS GIVEN THE OBSERVATIONS

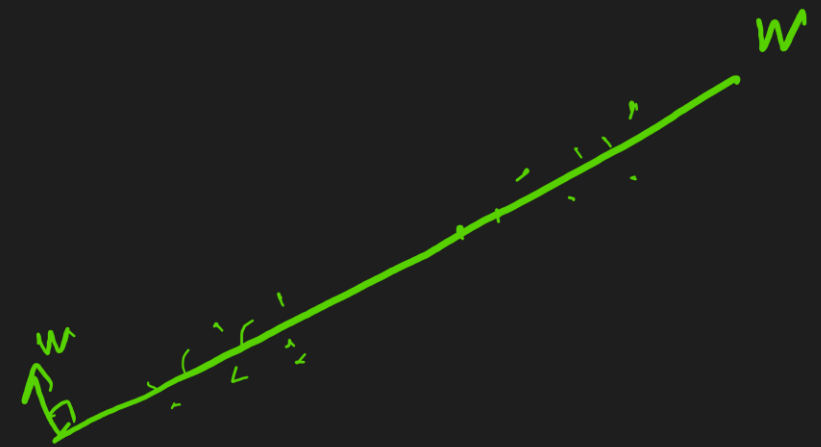
- $\alpha^{MAP} = \underset{\alpha}{\operatorname{argmax}} P(\alpha|X) = \underset{\alpha}{\operatorname{argmax}} \frac{P(X|\alpha) \cdot \underbrace{P(\alpha)}}{P(X)}$

- $\alpha^{MAP} = \underset{\alpha}{\operatorname{argmin}} \left(-\operatorname{LOG}(P(\alpha)) - \sum_{i=1}^n \operatorname{LOG} P(x^i|\alpha) \right)$

PRIOR VS POSTERIOR DISTRIBUTIONS

- $P(\alpha)$ CAPTURES THE **PRIOR** DISTRIBUTION
- $P(\alpha|X)$ CAPTURES THE **POSTERIOR** DISTRIBUTION
- IN OTHER WORDS,
 - WE START BY A PRIOR BELIEF ABOUT VALUE OF α
 - OUR BELIEF IS UPDATED AFTER SEEING SOME REAL DATA
 - THIS IS A BAYESIAN APPROACH

- $$P(y|x, W) = \left(\frac{1}{\gamma}\right) e^{-\frac{(x^T W - y)^2}{2\sigma^2}}$$



- $$P(x|W) = P(x)$$

$\sigma \ll 1$

- INPUT: $((x^1, y^1), \dots, (x^n, y^n))$ I.I.D. SAMPLE
 - EACH (x^i, y^i) IS DRAWN ACCORDING TO $P(X, Y)$
- MAXIMUM LIKELIHOOD ESTIMATE FOR W ?

ML ESTIMATE

$$\underset{w}{\operatorname{argmax}} p(Z|\alpha) = \underset{w}{\operatorname{argmax}} \sum \log p(x^i, y^i | w)$$

$$\approx \underset{w}{\operatorname{argmax}} \sum \log \left[\frac{p(x^i, y^i, w)}{p(w)} \right]$$

$$= \underset{w}{\operatorname{argmax}} \sum \log \left[p(y^i | x^i, w) p(x^i | w) \cancel{p(w)} / \cancel{p(w)} \right]$$

$$\approx \underset{w}{\operatorname{argmax}} \sum \left[\left(\ln \frac{1}{\sigma} \right) - \frac{(x^{iT} w - y^i)^2}{\sigma^2} + \ln p(x^i) \right]$$

$$= \underset{w}{\operatorname{argmin}} \sum_i (x^{iT} w - y^i)^2$$

A PROBABILISTIC MODEL FOR REGRESSION

- $P(y|x, W) = \left(\frac{1}{\gamma}\right) e^{-\frac{(x^T W - y)^2}{2\sigma^2}}$
- $P(x|W) = P(x)$
- MAXIMUM LIKELIHOOD SOLUTION FOR THIS PROBABILISTIC MODEL IS THE SAME AS LS SOLUTION.

MAXIMUM A POSTERIORI AND RLS

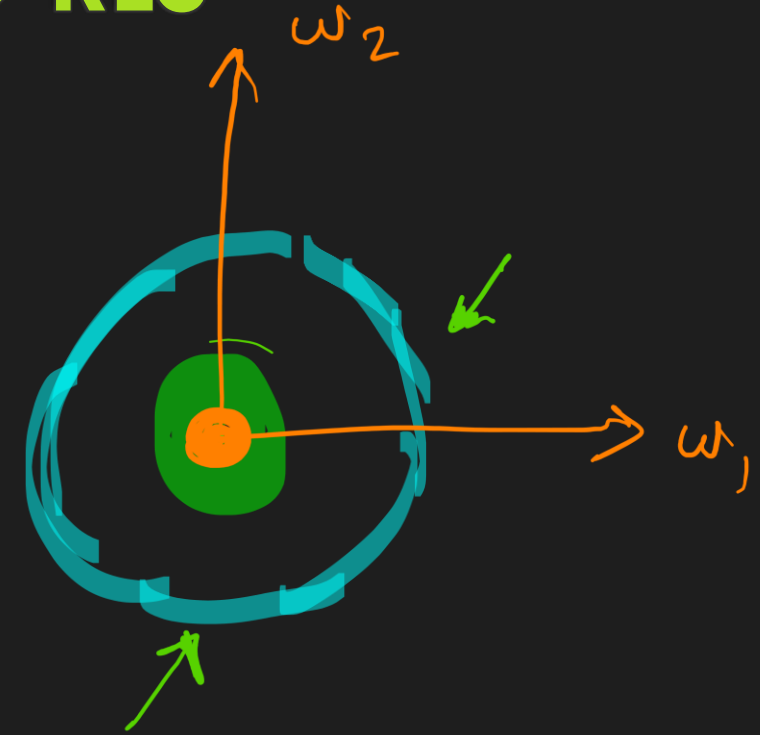
- $P(y|x, W) = \left(\frac{1}{\gamma}\right) e^{-\frac{(x^T W - y)^2}{2\sigma^2}}$

- $P(x|W) = P(x)$

- $P(W) = \frac{1}{\beta} e^{-\lambda \|W\|_2^2}$

what if we have

$$p(w) = \frac{1}{\beta} e^{-\frac{\lambda \|w\|_2^2}{\sigma^2}} \quad ?$$



MAP ESTIMATE

$$\operatorname{argmax}_w p(x|\alpha)P(\alpha) = \operatorname{argmax}_w \left[\underbrace{\log p(\alpha)} + \sum_i \underbrace{\log (x^i, y^i | w)} \right]$$

$$= \operatorname{argmax}_w \left[\underbrace{\ln \frac{1}{\sigma}} - \lambda \|w\|_2^2 - \frac{(x^{iT}w - y^i)^2}{\sigma^2} \right]$$

$$= \operatorname{argmin}_w \left[\sum_i (x^i w - y^i)^2 + \underbrace{\lambda \sigma^2}_{\downarrow} \|w\|_2^2 \right]$$

e.g., if $\sigma \geq 1$ then we get the standard RLS.

EXPECTED ERROR MINIMIZATION

- SQUARED LOSS

→ • $l(\hat{y}, y) = (y - \hat{y})^2$

- EMPIRICAL/TRAINING LOSS FOR $Z = \{(x^i, y^i)\}_{i=1}^n$

$$\frac{1}{n} \sum (y^i - \hat{y}^i)^2$$

- EXPECTED LOSS

$$L = \mathbf{E}_{\underbrace{(x,y)} \sim \underbrace{P(x,y)}} \underbrace{l(\hat{y}(x), y)} \leftarrow$$

REWRITING THE EXPECTATION

$$L = E_{(x,y) \sim p(x,y)} f(x,y) = E_{x \sim p(x)} E_{y \sim p(y|x)} f(x,y)$$

(Assuming the expectation exist)

Chain Rule

EXPECTED ERROR MINIMIZATION

- ASSUME WE KNOW $P(x, y)$

- $y^*(x) = \operatorname{argmin}_{\hat{y}} E_{x,y} (y - \hat{y}(x))^2 = ?$

$$\operatorname{argmin}_{\hat{y}(\cdot)} E_{x \sim P(x)} E_{y \sim P(y|x)} (y - \hat{y}(x))^2$$

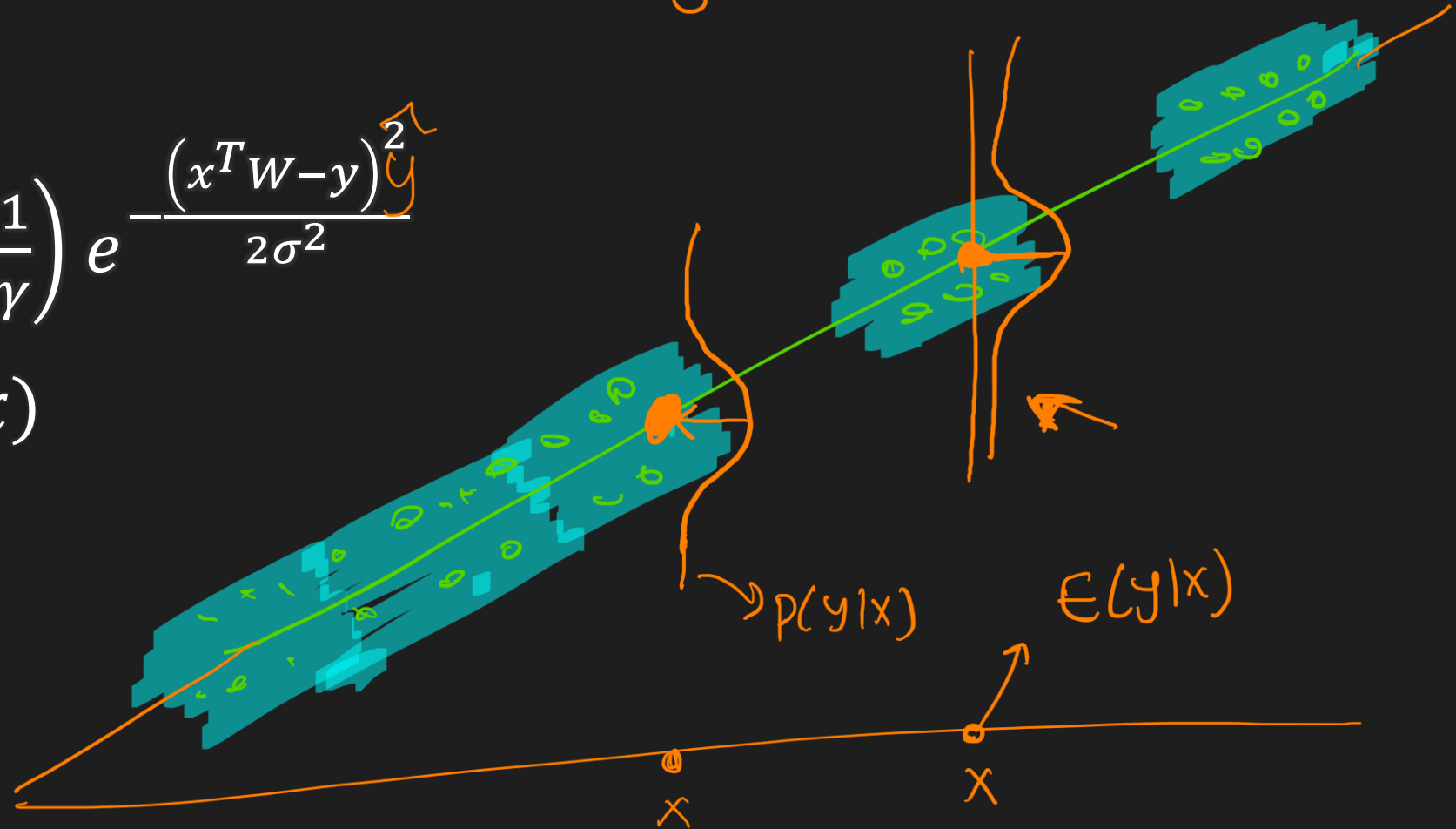
you can choose $\hat{y}(x)$ based on x

$$\rightarrow \operatorname{argmin}_{\substack{\hat{y}(x) \\ \in \mathbb{R}}} E_{y \sim P(y|x)} (y - \hat{y}(x))^2$$

EXAMPLE

- $P(y|x, W) = \left(\frac{1}{\gamma}\right) e^{-\frac{(x^T W - y)^2}{2\sigma^2}}$
- $P(x|W) = P(x)$

$$y^*(x) = E(y|x) = w^T x$$

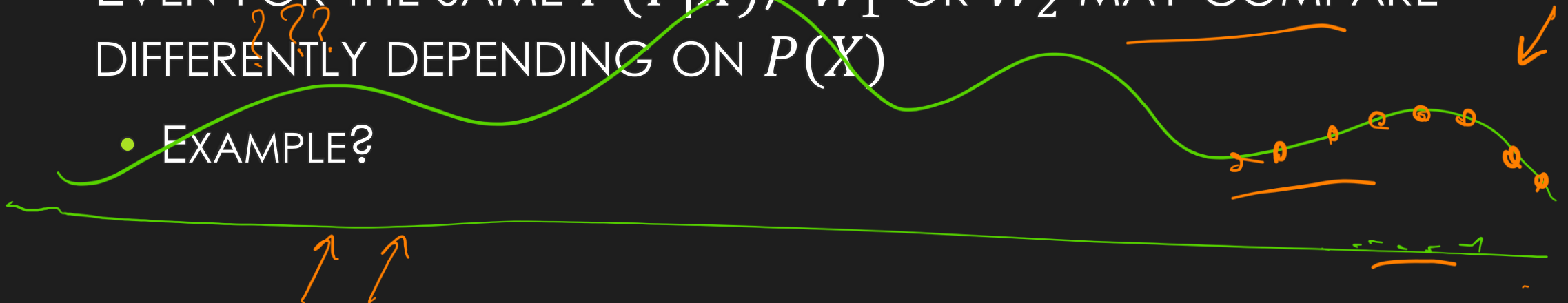


THE MARGINAL DISTRIBUTION MATTERS

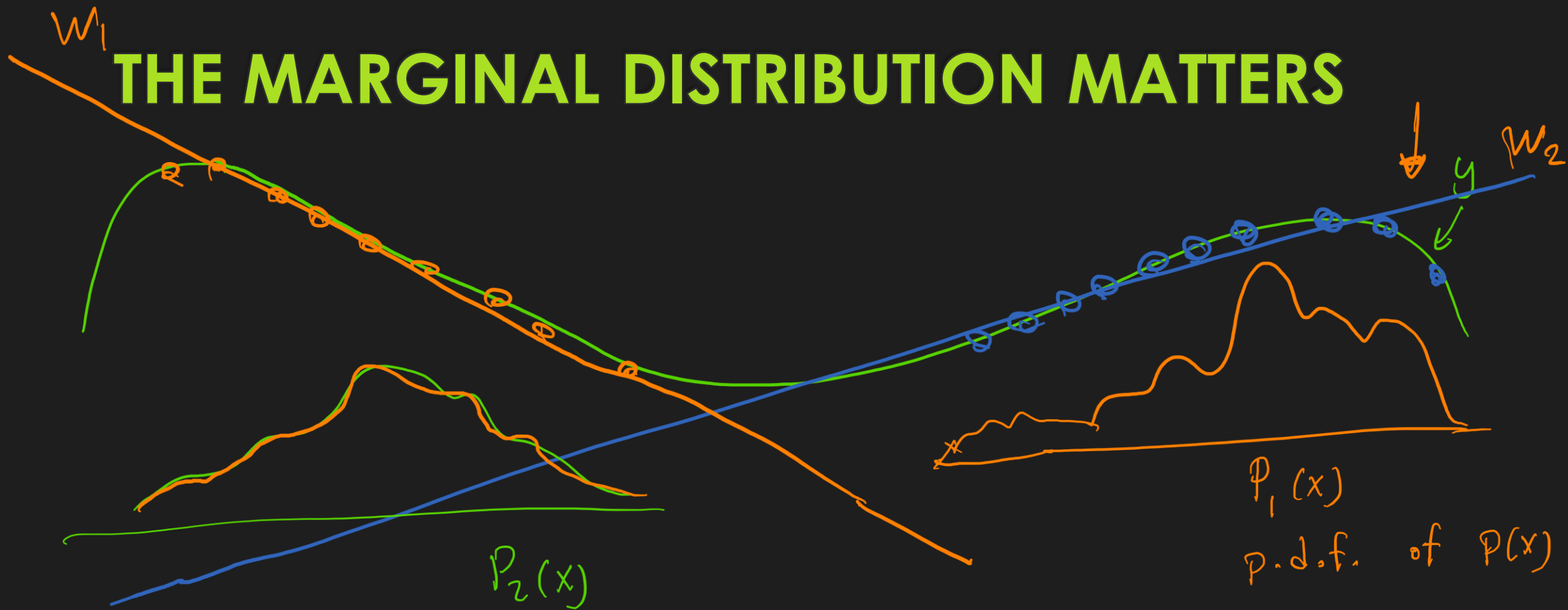
- IN PRACTICE, WE DON'T KNOW $P(X, Y)$
- ALSO, IT IS HARD TO ESTIMATE $E(Y|X)$ FOR EVERY X
 - BUT WE CAN AIM TO ESTIMATE $E(Y|X)$ "WHERE IT MATTERS"

- EVEN FOR THE SAME $P(Y|X)$, W_1 OR W_2 MAY COMPARE DIFFERENTLY DEPENDING ON $P(X)$

- EXAMPLE?



THE MARGINAL DISTRIBUTION MATTERS



Although despite $p(y|x)$ was the same.